



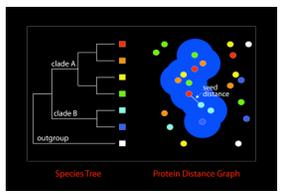
ABSTRACT

The ESPP2 project requires understanding of the microbial communities at contaminated field sites and, among other methods, will employ metagenomics in this endeavor. Metagenomics projects that seek to elucidate the population structure of microbial ecosystems are faced with the related computational challenges of classifying the sequences obtained and quantifying which organisms are present within a sample. Individually low-proportion species usually make up a large fraction of microbial communities, complicating their classification and quantification using traditional phylogenetic marker approaches. Such species usually don't yield sufficient read depth to assemble into longer sequences, leaving fragments that rarely contain traditional markers such as small subunit (SSU) rRNA gene. BLAST-based approaches for analysis of metagenomic sequences [1] compensate for this rarity of traditional markers, but may be confounded by genes that are subject to horizontal transfer or duplication. Another approach instead makes use of only reliable non-transferred single-copy genes [2] to classify and quantify the organisms present within a sample, but the application has so far been limited to the use of a fairly small set of universal genes found in all organisms. In this work, we have extended the latter approach, boosting the set of reliable marker genes from only about 30-40 universal genes to several hundred by identifying sets of single-copy genes that are not subject to inter-clade horizontal transfer through investigation of finished bacterial and archaeal genomes. These clade-oriented sequence markers allow for a method, which we have named "MicroCOSM", that greatly increases the probability that a marker will be found in any given sequence and therefore offers improved coverage for phylogenetic classification and quantification of microbial types in an environmental sample.

- [1] Huson D.H., Auch A.F., Qi J., Schuster S.C. (2007) "MEGAN analysis of metagenomic data." *Genome Res.* 17(3):377-86.
- [2] von Mering C., Hugenholz P., Raes J., Tringe S.G., Doerks T., Jensen L.J., Ward N., Bork P. (2007) "Quantitative phylogenetic assessment of microbial communities in diverse environments." *Science* 315(5815):1126-30.

DEVELOPMENT OF COSMS

Clade-oriented sequence marker (COSM) gene families are built by clustering of BLAST-detected homologous sequences. Clusters must not include genes that belong to species outside of the clade, determined by the largest nearest-neighbor distance between any two members of the cluster.



SPECIAL COSM TYPES: SCP AND NOVEL

While the COSMs may be used to assign membership to a clade, it is single-copy prevalent (SCP) genes that are best for branch placement within the clade, as they are not subject to duplication nor horizontal transfer. Additionally, novel families (NOVEL) that are only observed within a given clade may represent the introduction of protein families to that lineage.



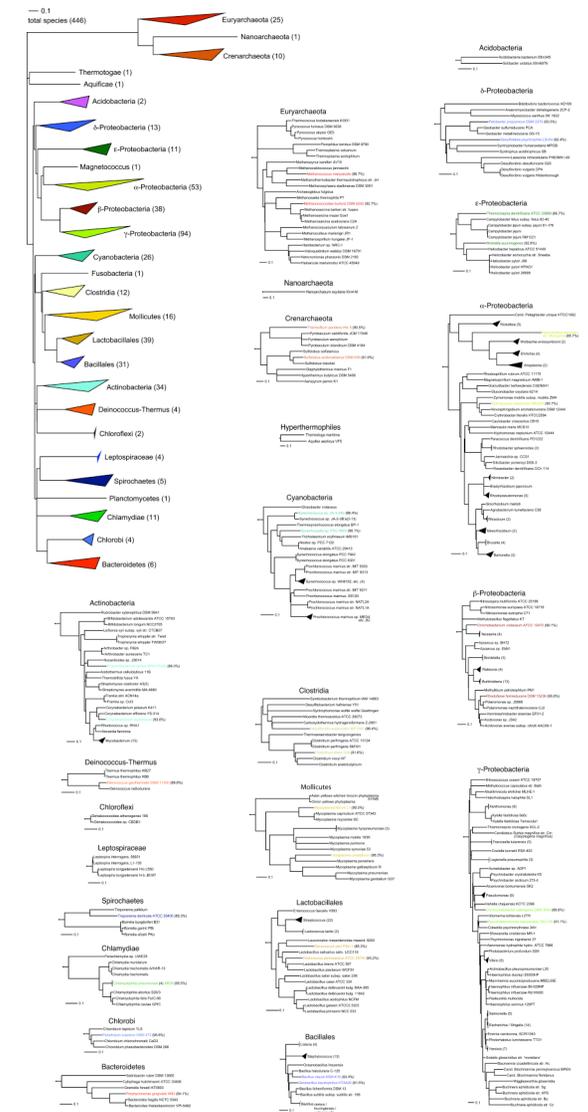
SCAN METAGENOMIC SEQUENCE WITH COSMS

In order to scan the metagenomic sequence with the clade-oriented sequence markers, profiles must be built and thresholds for membership in the protein family determined.

1. Make a multiple sequence alignment (MSA) of COSM proteins and build profile.
2. Determine lowest-scoring sequence for membership threshold.
3. Scan metagenomic reads with COSM profile and accept hits above threshold.
4. If COSM is also single copy prevalent (SCP), then fit into tree.

SPECIES TREE

Markers were developed from 146 species of Archaea and Bacteria with complete genomes. The 31 species that were removed for testing are shown in color.



COSM COUNTS

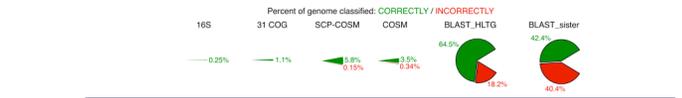
COSM, SCP and NOVEL counts for several high-level taxonomic groups (HTLG) specific to this data (not including children counts). Counts with the requirement of presence within >60% and >75% of clade member species are also reported.

High Level Taxonomic Group (HTLG)	COSM		SCP		Novel		High Level Taxonomic Group (HTLG)		COSM		SCP		Novel	
	>60%	>75%	>60%	>75%	present	present	present	present	>60%	>75%	present	present	>60%	>75%
Bacteria + Archaea	N/A	N/A	54	37	N/A	N/A	N/A	N/A	224	24	11	8	1	4
Archaea	1218	206	135	152	101	308	66	55	** Firmicutes	858	12	13	4	285
Euryarchaeota	328	23	17	16	11	34	7	7	** Basili	8	0	0	0	0
Crenarchaeota	400	58	32	36	19	109	15	11	** Bacilli	224	24	11	8	1
Bacteria	738	184	134	83	63	180	40	32	** Lactobacillales	332	21	8	7	97
Hyperthermophiles	27	27	27	26	26	1	1	1	** Bacillales	172	22	15	13	11
Proteobacteria + Asso	1205	207	144	0	0	195	1	0	** Supergroup 1	207	200	193	152	37
Actinobacteria	397	397	397	317	317	103	103	103	** Chiroflexi	366	366	340	340	112
Proteobacteria	2283	35	11	0	957	2	0	0	** Bacteroidetes	186	4	1	4	33
** Delta	211	8	6	2	27	1	1	1	** Supergroup 1 + Lepto	159	13	8	4	3
** Epsilon	306	144	118	13	2	101	40	35	** Leptospiraceae	1159	1137	1028	1013	498
** Alpha	114	0	24	0	28	6	4	4	** Chlamydiae	159	13	8	4	3
** Beta + Gamma	247	28	15	2	895	4	0	0	** Chlamydiae	174	174	166	159	52
** Beta	151	24	19	0	35	4	3	3	** Bacteroidetes-Chlorobi	100	34	27	20	18
** Gamma	421	13	9	10	7	121	3	2	** Chlorobi	262	252	180	203	149
Cyanobacteria	522	140	122	108	96	250	80	68	** Bacteroidetes	92	41	30	19	22

COVERAGE AND ACCURACY

We removed 31 species that were of varying distance by 16S rRNA gene similarity to the closest remaining species and built a test set of COSMs. We then examined the fraction of each of these 31 genomes that could accurately be assigned to high-level taxonomic groups (e.g. phyla) using 16S rRNA genes, the 31 COGs used by von Mering *et al.*, our test COSMs and SCP-COSMs (both with >75% prevalence requirement), and a BLAST comparison with the proteomes of the remaining species. We also examined more fine-grained taxonomic assignment at the sister level with BLAST as the only method roughly comparable to phylogenetic placement.

Species	HTLG	Nearest Relative (IS%)	Genome Size (Mb)	16S COV(%)	SCP COV(%)	SCP HTLG	SCP HTLG	COSM COV(%)	COSM HTLG	BLAST COV(%)	BLAST HTLG	BLAST Sister	BLAST Accn(%)	BLAST Accn(%)
Methanococcus marisnigri	Euryarchaeota	86.7	1.66	0.25	1.5	6.6	100	4.3	95.1	85.6	87.5	85.6	85.6	85.6
Methanococcus burtoni DSM 6242	Euryarchaeota	92.7	2.58	0.17	1.4	5.4	98.2	6	88.9	77	88.8	77	80.1	80.1
Thermotoga parvula H8.5	Crenarchaeota	95.3	1.81	0.083	1.4	5.1	99.2	2.6	92.5	75.5	88.8	75.5	81.4	81.4
Sulfolobus acidocaldarius DSM 639	Crenarchaeota	93.8	2.23	0.066	1.2	14.9	99.3	14.8	96.7	81.5	84.8	81.5	66.1	66.1
Desulfotribes psychrophilus L3454	Delta	85.4	2.66	0.45	0.59	3.3	93.1	1.9	61.5	75.9	82.3	75.9	17.5	17.5
Pelobacter propionicus DSM 2379	Delta	95.5	4.24	0.15	0.53	4.2	95.7	2.7	75.1	82.9	86.9	82.9	55.9	55.9
Thermoplasma denitrificans ATCC 33888	Epsilon	86.7	2.2	0.29	2.1	4.2	100	1.5	81.4	84.4	85.8	84.4	40	40
Wohlfelia succrogenes	Epsilon	92.8	2.11	0.42	1.1	3.7	100	1.1	88	90.9	84.9	90.9	43.9	43.9
Spirochaeta aurantiflava PR2528	Alpha	85.7	3.86	0.17	2.5	2.7	100	1.5	100	75	81.5	75	72.8	72.8
Chromohalobacter volcanodurans ATCC 12472 B1	Gamma	90.1	4.75	0.25	0.48	2.2	94.8	0.7	59	83	86.5	83	11.4	11.4
Chromohalobacter volcanodurans DSM 3043	Gamma	86.6	3.7	0.2	0.61	3.1	97.2	1.2	74.7	87	87.2	87	10.6	10.6
Planctomycetes halophilus TAC125	Gamma	90.1	3.85	0.15	0.58	3.3	93.3	1.2	77.7	84	87.4	84	24.5	24.5
Syntherobacterales sp. PCC 6803	Cyanobacteria	91.1	3.95	0.15	0.61	4.1	98.3	3.3	93.6	81.6	83.4	81.6	62.3	62.3
Nitrospira carolinensis ATCC 37426	Alpha	85.7	3.86	0.17	2.5	2.7	100	1.5	100	75	81.5	75	72.8	72.8
Desulfotribes psychrophilus MP104C	Clostridia	86.4	2.35	0.14	0.95	3.95	95.5	4.4	80.7	85.3	80.7	85.3	17.2	17.2
Chlorobium thauri DSM	Mollicutes	91.6	2.67	0.63	0.74	3.5	94.4	1.2	67.4	81	82.1	81	38.2	38.2
Ureaplasma urealyticum	Mollicutes	86.3	0.75	0.79	3.9	97.1	0.25	29.8	80.7	83.4	80.7	83.4	33.9	33.9
Mycoplasma thermophilum 1.1	Mollicutes	96	0.76	0.19	2.3	9.6	98.3	3.8	97.7	86.6	83.2	86.6	37.6	37.6
Denitrosospora sp. PSJ-1	Lactobacillales	85.3	1.78	0.18	1.3	5	100	1.3	88	79	91.1	79	46.1	46.1
Nitrospira carolinensis ATCC 37426	Lactobacillales	85.7	0.83	0.13	0.13	1.2	100	1.2	80	84.9	82	84.9	59.4	59.4
Syntherobacterales sp. PCC 6803	Bacillales	91.5	3.99	0.39	0.51	96.1	3.5	80.7	80.9	77.3	80.9	77.3	59.4	59.4
Bacillus clausii DSM 145	Bacillales	85.8	4.3	0.25	0.56	5.7	97	1.9	87	88.8	83.8	81	47.5	47.5
Prophococcus aureus KPA17102	Actinobacteria	82.3	2.56	0.18	0.91	5.1	98.4	2.4	96.7	82.8	81.4	82.8	28.8	28.8
Coniobacterium glomeratum	Actinobacteria	83.8	2.51	0.26	0.84	7.2	99.2	3.1	87.8	83.8	83.8	81	83.8	83.8
Denitrosospora sp. PCC 6803	Denitrosospora	89.8	3.04	0.2	0.76	7.2	97.5	5.7	88.6	87	78.6	87	73.5	73.5
Topoglossina aerifera ATCC 3640	Spirochaetes	85.6	2.36	0.12	0.86	3.8	97.4	1.1	80.4	77.1	80.7	77.1	30.7	30.7
Chlamydia pneumoniae AF99	Chlamydiae	95.5	1.23	0.12	1.9	26.1	100	22.7	100	82.3	87.6	82.3	96.3	96.3
Pelotribion laevis DSM 223	Citrodia	85.6	2.36	0.12	0.86	3.8	97.4	1.1	80.4	77.1	80.7	77.1	30.7	30.7
Porphyromonas gingivalis W83	Bacteroidetes	84.1	2.34	0.25	0.96	3.5	93.7	5.8	61.1	76	84	76	77	77
AVERAGE		90.7	2.71	0.25	1.1	5.7	96.5	3.8	88	82.7	84.9	82.7	51.2	51.2



CONCLUSIONS

We have improved the coverage provided by traditional universal sequence markers by identifying clade-oriented sequence markers. Using stringent requirements for inclusion in the gene family, coverage increases from the ~1% of universal COGs to approximately 3-6% with COSMs, accompanied by a small loss in accuracy. Greater coverage but less accuracy is expected using looser thresholds. Comparison with pair-wise BLAST-based methods shows COSMs to be far more accurate, albeit with less coverage. Unlike, unlike BLAST-based methods, SCP-COSMs permit placement explicitly on the tree, and permit greater coverage with which to assess the population structure of a microbial community than universal COGs alone. We expect that more comprehensive sequencing in poorly sampled clades, such as Clostridia, will improve results in future versions of the method.

ACKNOWLEDGEMENTS

supported by the U. S. Department of Energy, Office of Science, Office of Biological and Environmental Research, Genomics:GTL Program through contract DE-AC02-05CH11231 between Lawrence Berkeley National Laboratory and the U. S. Department of Energy.